

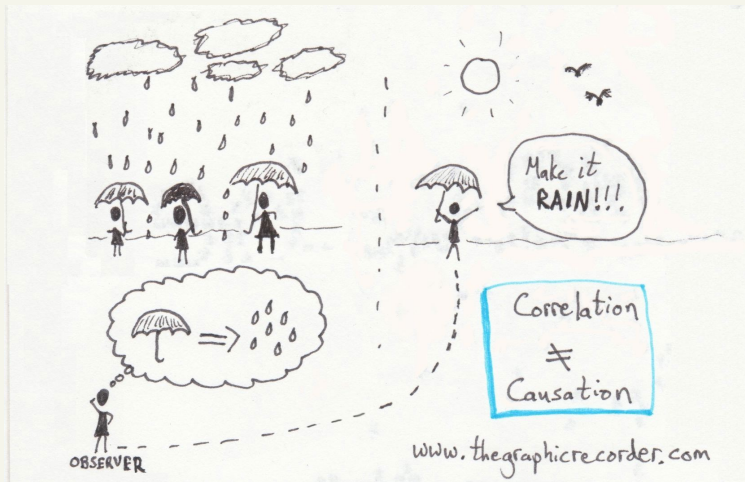
Modelos Bayesianos

Introducción a la Inteligencia Artificial

Agenda

- 1 Inferencia Bayesiana
- 2 Naive Bayes y Gaussian Naive Bayes
- 3 Multinomial Naive Bayes
- 4 Muestreo y PyMC

Hablemos de causalidad



Causalidad: ¿Quién causa a quién?

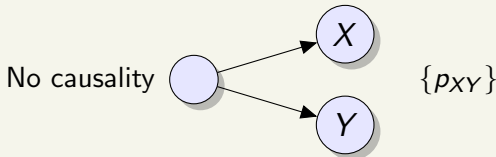
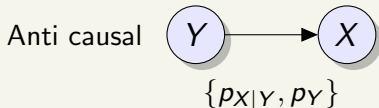
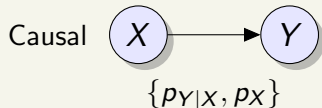
Independent Causal Mechanisms (ICM) Principle

The causal generative process of a system's variables is composed of autonomous modules that do not inform or influence each other. In the probabilistic case, this means that the conditional distribution of each variable given its causes (i.e., its mechanism) does not inform or influence the other mechanisms.

Causalidad: ¿Quién causa a quién?

Independent Causal Mechanisms (ICM) Principle

The causal generative process of a system's variables is composed of autonomous modules that do not inform or influence each other. In the probabilistic case, this means that the conditional distribution of each variable given its causes (i.e., its mechanism) does not inform or influence the other mechanisms.



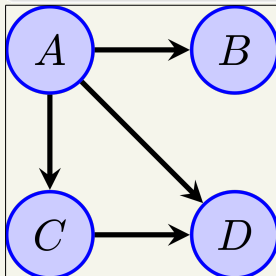
Redes Bayesianas

Modelos Gráficos

Modelos probabilísticos capaz de representarse con un grafo.

Red Bayesiana

Grafo acíclico dirigido que representa la relación de causalidad e independencia de sus variables. Dos variables aleatorias cualesquiera son condicionalmente independientes dados los valores de sus padres causales (y por lo tanto las raíces son independientes).



$$p(A, B, C, D) = p(A) \cdot p(B|A) \cdot p(C|A) \cdot p(D|A, C)$$

Inferencia Bayesiana

- Los parámetros θ deben ser considerado realizaciones de una variable aleatoria T con una distribución a priori conocida $p(\theta)$.
- Las muestras son i.i.d. **cuando** se conoce el parámetro.
- La distribución a posteriori de los parámetros se calcula como:

$$p(\theta|\mathcal{D}_n) \propto p(\theta) \prod_{i=1}^n p(x_i|\theta)$$

con $\mathcal{D}_n = \{x_1, \dots, x_n\}$.

- Como estimador puntual suele elegirse el *maximo a posteriori* (Θ discreto) y el estimador bayesiano o *media a posteriori* (Θ continuo).
- No son necesarios los estimadores puntuales para predecir:

$$p(x_{\text{test}}|\mathcal{D}_n) = \int_{\Theta} p(x_{\text{test}}|\theta)p(\theta|\mathcal{D}_n)d\theta = \mathbb{E}[p(x_{\text{test}}|T)|\mathcal{D}_n]$$

Inferencia Bayesiana

Filosofía Bayesiana

La estadística bayesiana interpreta la probabilidad como una medida de credibilidad en un evento. Por eso se habla de que el enfoque Bayesiano busca verdades en contexto de incertidumbre.

Inferencia Bayesiana

Filosofía Bayesiana

La estadística bayesiana interpreta la probabilidad como una medida de credibilidad en un evento. Por eso se habla de que el enfoque Bayesiano busca verdades en contexto de incertidumbre.

No confundir Bayesiano con Relativista

¿Es posible entonces alcanzar verdades en las ciencias empíricas en las que es inevitable decir “no sé”? Sí. Podemos evitar mentir: maximizando incertidumbre (no afirmar más de lo que se sabe) dada la información disponible (sin ocultar lo que sí se sabe).

Inferencia Bayesiana

Filosofía Bayesiana

La estadística bayesiana interpreta la probabilidad como una medida de credibilidad en un evento. Por eso se habla de que el enfoque Bayesiano busca verdades en contexto de incertidumbre.

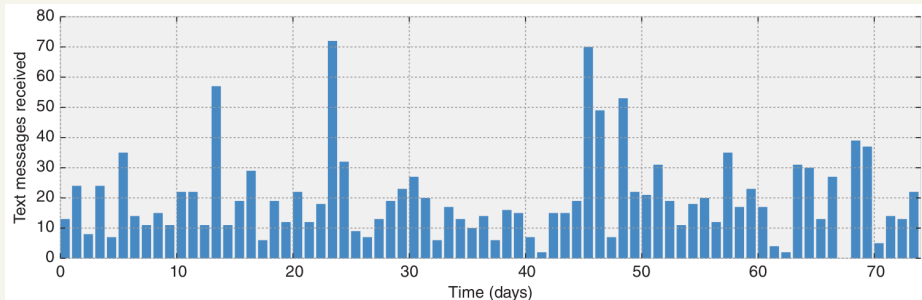
No confundir Bayesiano con Relativista

¿Es posible entonces alcanzar verdades en las ciencias empíricas en las que es inevitable decir “no sé”? Sí. Podemos evitar mentir: maximizando incertidumbre (no afirmar más de lo que se sabe) dada la información disponible (sin ocultar lo que sí se sabe).

¿Son prácticos los métodos Bayesianos?

Si, no solo por poder *adaptarse* a intentar resolver los mismos problemas que la estadística frecuentista (por ejemplo predicciones), sino que también pueden intentar resolver problemas donde la estadística clásica es insuficiente o iluminar el sistema subyacente con un modelado más flexible.

Ejemplo de modelado Bayesiano



Problema

Un usuario proporciona una serie de recuentos diarios de mensajes de whatsapp enviados. Tiene curiosidad por saber si los hábitos de envío de mensajes han cambiado con el tiempo. ¿Cómo puedes modelar esto?

Ejemplo de modelado Bayesiano

- La cantidad de mensajes en un día deberá ser modelada como una variable discreta cuyos átomos es \mathbb{N}_0 . Por ejemplo $X_i \sim \text{Poi}(\lambda_i)$.
- Si observamos los datos, parecería que el valor de λ_i aumenta en algún momento durante las observaciones. ¿Cómo podemos representar matemáticamente esta observación? Supongamos que algún día τ durante el período de observación, el parámetro λ_i se incrementa repentinamente. Entonces realmente tenemos dos tasas: una para el período anterior a τ y otro para el resto del período:

$$\lambda_i = \begin{cases} \beta_1 & i < \tau \\ \beta_2 & i \geq \tau \end{cases}$$

- Tanto β_1 como β_2 toman valores reales no negativos. Por ejemplo $\beta_1, \beta_2 \sim \mathcal{E}(\alpha)$. Nuestra estimación de α no influye demasiado en el modelo, por lo que tenemos cierta flexibilidad en nuestra elección. Para evitar ser demasiado obstinados con este parámetro se sugiere:

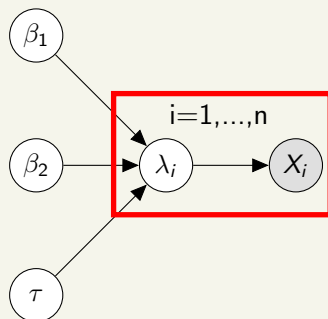
$$\alpha \approx \frac{1}{\frac{1}{n} \sum_{i=1}^n X_i}$$

Ejemplo de modelado Bayesiano

- ¿Qué pasa con τ ? Debido a la varianza de los datos, es difícil caracterizarlo en detalle. En cambio, podemos asignar la creencia menos informativa posibles $\tau \sim \mathcal{U}\{1 : T\}$.

Ejemplo de modelado Bayesiano

- ¿Qué pasa con τ ? Debido a la varianza de los datos, es difícil caracterizarlo en detalle. En cambio, podemos asignar la creencia menos informativa posibles $\tau \sim \mathcal{U}\{1 : T\}$.



Outline

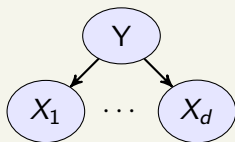
- 1 Inferencia Bayesiana
- 2 Naive Bayes y Gaussian Naive Bayes**
- 3 Multinomial Naive Bayes
- 4 Muestreo y PyMC

Naive Bayes

Naive Bayes

Estimar los parámetros (máxima verosimilitud o bayesiano) asumiendo que una relación de causalidad $Y \rightarrow X$ con las diferentes componentes $X_j | Y = k$ independientes.

Red Bayesiana



Cálculo

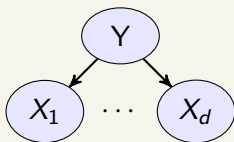
$$p(y|x) = \frac{p(y) \prod_{j=1}^d p(x_j|y)}{p(x)}$$

Naive Bayes

Naive Bayes

Estimar los parámetros (máxima verosimilitud o bayesiano) asumiendo que una relación de causalidad $Y \rightarrow X$ con las diferentes componentes $X_j|Y = k$ independientes.

Red Bayesiana



Cálculo

$$p(y|x) = \frac{p(y) \prod_{j=1}^d p(x_j|y)}{p(x)}$$

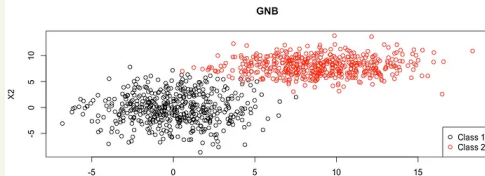
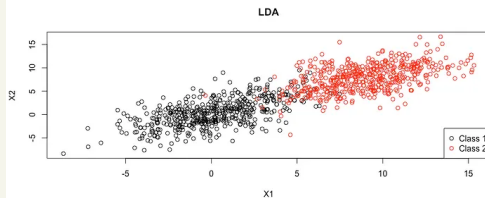
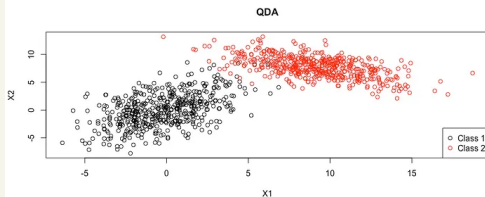
Gaussian Naive Bayes

Sea $\pi = \{c_1, \dots, c_K\}$ y $\theta^{(k)} = \{(\mu_1^{(k)}, \sigma_1^{2(k)}), \dots, (\mu_d^{(k)}, \sigma_d^{2(k)})\}$, se modela $Y \sim \text{Cat}(\{c_1, \dots, c_K\})$ y $X_j|Y = k \sim \mathcal{N}(\mu_j^{(k)}, \sigma_j^{2(k)})$, para luego estimar los parámetros.

Gaussian Naive Bayes (GNB)

Diferencias entre QDA, LDA y GNB

- QDA acepta como Σ_k cualquier conjunto de matrices definidas positiva.
- LDA acepta como Σ_k cualquier matriz pero todas iguales.
- GNB permite tener matrices Σ_k diferentes pero todas diagonales.



Gaussian Naive Bayes (GNB)

QDA

$$\Sigma_k = \frac{1}{|\mathcal{D}_k| - 1} \sum_{x \in \mathcal{D}_k} (x - \mu^{(k)}) (x - \mu^{(k)})^T$$

LDA

$$\Sigma = \frac{1}{n - K} \sum_{k=1}^K (|\mathcal{D}_k| - 1) \Sigma_k$$

GNB

$$\Sigma_k = \text{DIAG} \left(\sigma_1^{2(k)}, \dots, \sigma_d^{2(k)} \right), \quad \sigma_j^{2(k)} = \frac{1}{|\mathcal{D}_k| - 1} \sum_{x \in \mathcal{D}_k} (x_j - \mu_j^{(k)})^2$$

Outline

- 1 Inferencia Bayesiana
- 2 Naive Bayes y Gaussian Naive Bayes
- 3 Multinomial Naive Bayes**
- 4 Muestreo y PyMC

Multinomial Naive Bayes

Multinomial Naive Bayes (MNB)

Sea $\pi = \{c_1, \dots, c_K\}$ y $\theta^{(k)} = \{\theta_1^{(k)}, \dots, \theta_V^{(k)}\}$, se modela como $Y \sim \text{Cat}(\{c_1, \dots, c_K\})$ y $X_j | Y = k \sim \text{Cat}(\{\theta_1^{(k)}, \dots, \theta_V^{(k)}\})$, utilizando estimadores puntuales.

Multinomial Naive Bayes

Multinomial Naive Bayes (MNB)

Sea $\pi = \{c_1, \dots, c_K\}$ y $\theta^{(k)} = \{\theta_1^{(k)}, \dots, \theta_V^{(k)}\}$, se modela como $Y \sim \text{Cat}(\{c_1, \dots, c_K\})$ y $X_j | Y = k \sim \text{Cat}(\{\theta_1^{(k)}, \dots, \theta_V^{(k)}\})$, utilizando estimadores puntuales.

Inferencia

$$p(y|x) \propto c_y \prod_{j=1}^d \theta_{x_j}^{(y)} = c_y \prod_{m=1}^V \left(\theta_m^{(y)} \right)^{N_m}$$

con N_m : Cantidad de predictores con valor m .

Multinomial Naive Bayes

Multinomial Naive Bayes (MNB)

Sea $\pi = \{c_1, \dots, c_K\}$ y $\theta^{(k)} = \{\theta_1^{(k)}, \dots, \theta_V^{(k)}\}$, se modela como $Y \sim \text{Cat}(\{c_1, \dots, c_K\})$ y $X_j|Y = k \sim \text{Cat}(\{\theta_1^{(k)}, \dots, \theta_V^{(k)}\})$, utilizando estimadores puntuales.

Inferencia

$$p(y|x) \propto c_y \prod_{j=1}^d \theta_{x_j}^{(y)} = c_y \prod_{m=1}^V \left(\theta_m^{(y)} \right)^{N_m}$$

con N_m : Cantidad de predictores con valor m .

Sobre las variables contadoras

Sea $N = (N_1, \dots, N_V)$, es sencillo notar que $\sum_{m=1}^V N_m = d$ y $N|_{Y=k} \sim \mathcal{M}_n(d, [\theta_1^{(k)}, \dots, \theta_V^{(k)}])$.

Multinomial Naive Bayes

En inferencia se elige el máximo de una transformación afín

$$\arg \max_y p(y|x) = \arg \max_y \log(c_y) + \sum_{m=1}^V N_m \log(\theta_m^{(y)})$$

Multinomial Naive Bayes

En inferencia se elige el máximo de una transformación afín

$$\arg \max_y p(y|x) = \arg \max_y \log(c_y) + \sum_{m=1}^V N_m \log(\theta_m^{(y)})$$

Probabilidades de las Clases

Los parámetros c_1, \dots, c_K son estimados por máxima verosimilitud como:

$$\hat{c}_k = \frac{\#\{y_i = k\}}{n}$$

Multinomial Naive Bayes

En inferencia se elige el máximo de una transformación afín

$$\arg \max_y p(y|x) = \arg \max_y \log(c_y) + \sum_{m=1}^V N_m \log(\theta_m^{(y)})$$

Probabilidades de las Clases

Los parámetros c_1, \dots, c_K son estimados por máxima verosimilitud como:

$$\hat{c}_k = \frac{\#\{y_i = k\}}{n}$$

Sobre el valor d

Cada muestra puede poseer un valor de d diferente. Eso es típico en texto, donde cada documento posee una cantidad diferente de palabras.

Multinomial Naive Bayes

Estimación de $\theta_m^{(k)}$

Se cuenta con datos $\{(N_i, y_i)\}_{i=1}^n$. Sin embargo, para cada clase k se utilizarán solamente los datos con $\{y_i = k\}$ distribuidos como una multinomial de probabilidades $\theta_1^{(k)}, \dots, \theta_V^{(k)}$. A su vez, dado que las variables N_m cuentan ocurrencias, puedo compactar todas las muestras de entrenamiento de cada clase en una sola (suficiencia estadística).

$$\tilde{N}_m^{(k)} = \sum_{i=1}^n N_{i,m} \cdot \mathbb{1}\{y_i = k\}$$

Multinomial Naive Bayes

Estimación de $\theta_m^{(k)}$

Se cuenta con datos $\{(N_i, y_i)\}_{i=1}^n$. Sin embargo, para cada clase k se utilizarán solamente los datos con $\{y_i = k\}$ distribuidos como una multinomial de probabilidades $\theta_1^{(k)}, \dots, \theta_V^{(k)}$. A su vez, dado que las variables N_m cuentan ocurrencias, puedo compactar todas las muestras de entrenamiento de cada clase en una sola (suficiencia estadística).

$$\tilde{N}_m^{(k)} = \sum_{i=1}^n N_{i,m} \cdot \mathbb{1}\{y_i = k\}$$

Modelado: Estimador Bayesiano

Como modelado para el entrenamiento se supone $T \sim \text{Dir}([\alpha_1, \dots, \alpha_V])$ y $(\tilde{N}_1^{(k)}, \dots, \tilde{N}_V^{(k)})|_{T=\theta} \sim \mathcal{M}_n(\tilde{d}^{(k)}, [\theta_1^{(k)}, \dots, \theta_V^{(k)}])$.

Multinomial Naive Bayes

Dirichlett Distribution

El vector aleatorio $(T_1, \dots, T_V) \sim \text{Dir}([\alpha_1, \dots, \alpha_V])$ puede ser pensado como una beta multivariada. Su densidad es de la forma

$$p(\theta_1, \dots, \theta_V) = \frac{\prod_{m=1}^V \Gamma(\alpha_m)}{\Gamma\left(\sum_{m=1}^V \alpha_m\right)} \left(\prod_{m=1}^V \theta_m^{\alpha_m-1} \right) \cdot \mathbb{1} \left\{ \sum_{m=1}^V \theta_m = 1, \theta_m \geq 0 \right\}$$

con sus marginales $T_m \sim \beta(\alpha_m, \sum_{\eta \neq m} \alpha_\eta)$.

Sobre la beta

Recordar que si $T \sim \beta(a, b)$, entonces $\mathbb{E}[T] = \frac{a}{a+b}$.

Multinomial Naive Bayes

Distribución a Posteriori

$$\begin{aligned} & p\left(\theta_1^{(k)}, \dots, \theta_V^{(k)} \mid \tilde{N}_1^{(k)}, \dots, \tilde{N}_V^{(k)}\right) \\ & \propto P\left(\tilde{N}_1^{(k)}, \dots, \tilde{N}_V^{(k)} \mid \theta_1^{(k)}, \dots, \theta_V^{(k)}\right) \cdot p\left(\theta_1^{(k)}, \dots, \theta_V^{(k)}\right) \\ & \propto \left(\prod_{m=1}^V (\theta_m^{(k)})^{\tilde{N}_m^{(k)}}\right) \left(\prod_{m=1}^V (\theta_m^{(k)})^{\alpha_m - 1} \cdot \mathbb{1}\left\{\theta_m^{(k)} \geq 0\right\}\right) \cdot \mathbb{1}\left\{\sum_{m=1}^V \theta_m^{(k)} = 1\right\} \end{aligned}$$

con lo cual $\mathbf{T} \mid \tilde{N}_1^{(k)}, \dots, \tilde{N}_V^{(k)} \sim \text{Dir}([\tilde{N}_1^{(k)} + \alpha_1, \dots, \tilde{N}_V^{(k)} + \alpha_V])$

Multinomial Naive Bayes

Distribución a Posteriori

$$\begin{aligned} p\left(\theta_1^{(k)}, \dots, \theta_V^{(k)} \mid \tilde{N}_1^{(k)}, \dots, \tilde{N}_V^{(k)}\right) \\ \propto P\left(\tilde{N}_1^{(k)}, \dots, \tilde{N}_V^{(k)} \mid \theta_1^{(k)}, \dots, \theta_V^{(k)}\right) \cdot p\left(\theta_1^{(k)}, \dots, \theta_V^{(k)}\right) \\ \propto \left(\prod_{m=1}^V (\theta_m^{(k)})^{\tilde{N}_m^{(k)}}\right) \left(\prod_{m=1}^V (\theta_m^{(k)})^{\alpha_m - 1} \cdot \mathbb{1}\left\{\theta_m^{(k)} \geq 0\right\}\right) \cdot \mathbb{1}\left\{\sum_{m=1}^V \theta_m^{(k)} = 1\right\} \end{aligned}$$

con lo cual $T \mid \tilde{N}_1^{(k)}, \dots, \tilde{N}_V^{(k)} \sim \text{Dir}([\tilde{N}_1^{(k)} + \alpha_1, \dots, \tilde{N}_V^{(k)} + \alpha_V])$

Estimador Bayesiano

$$\hat{\theta}_m^{(k)} = \mathbb{E}[T_m \mid \tilde{N}_1^{(k)}, \dots, \tilde{N}_V^{(k)}] = \frac{\tilde{N}_m^{(k)} + \alpha_m}{\sum_{\eta=1}^V \tilde{N}_\eta^{(k)} + \alpha_\eta}$$

Outline

- 1 Inferencia Bayesiana
- 2 Naive Bayes y Gaussian Naive Bayes
- 3 Multinomial Naive Bayes
- 4 Muestreo y PyMC

Monte-Carlo

Ley de los grandes números

Sean X_i variables aleatorias con esperanza finita $\mathbb{E}[X]$, entonces $\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{n \rightarrow \infty} \mathbb{E}[X]$ (w.p.1).

Monte-Carlo

Ley de los grandes números

Sean X_i variables aleatorias con esperanza finita $\mathbb{E}[X]$, entonces $\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{n \rightarrow \infty} \mathbb{E}[X]$ (w.p.1).

Método de Monte-Carlo

Sean X_i variables aleatorias i.i.d con pdf $p(x)$ o pmf $P(x)$, entonces

$$\int_{\mathbb{R}} g(x)p(x)dx \approx \frac{1}{n} \sum_{i=1}^n g(X_i), \quad \sum_{x \in \mathbb{A}} g(x)P(x) \approx \frac{1}{n} \sum_{i=1}^n g(X_i)$$

Monte-Carlo

Ley de los grandes números

Sean X_i variables aleatorias con esperanza finita $\mathbb{E}[X]$, entonces $\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{n \rightarrow \infty} \mathbb{E}[X]$ (w.p.1).

Método de Monte-Carlo

Sean X_i variables aleatorias i.i.d con pdf $p(x)$ o pmf $P(x)$, entonces

$$\int_{\mathbb{R}} g(x)p(x)dx \approx \frac{1}{n} \sum_{i=1}^n g(X_i), \quad \sum_{x \in \mathbb{A}} g(x)P(x) \approx \frac{1}{n} \sum_{i=1}^n g(X_i)$$

Algunas variantes interesantes

- $\int_a^b g(x)dx \approx \frac{b-a}{n} \sum_{i=1}^n g(X_i)$ con $X_i \sim \mathcal{U}(a, b)$.
- $\int_a^b g(x)dx \approx \frac{1}{kn} \sum_{i=1}^n \mathbb{1}\{a < X_i < b\}$ con X_i de pdf $p(x) = k \cdot g(x)$.
- $\int_{\mathbb{R}} g(x)p(x)dx \approx \frac{1}{n} \sum_{i=1}^n \frac{p(X_i)}{q(X_i)} g(X_i)$ con X_i de pdf $q(x)$.

Técnicas de Muestreo

Ideas Principales

- En lugar de obtener la distribución a posteriori, vamos a generar muestras de esta distribución.
- Vamos aproximar la predictiva por Monte-Carlo.

$$p(x_{\text{test}}|\mathcal{D}_n) \approx \frac{1}{N_s} \sum_{i=1}^{N_s} p(x_{\text{test}}|T_i)$$

Técnicas de Muestreo

Ideas Principales

- En lugar de obtener la distribución a posteriori, vamos a generar muestras de esta distribución.
- Vamos aproximar la predictiva por Monte-Carlo.

$$p(x_{\text{test}}|\mathcal{D}_n) \approx \frac{1}{N_s} \sum_{i=1}^{N_s} p(x_{\text{test}}|T_i)$$

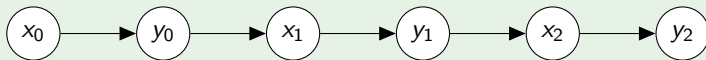
Muestreo de Gibbs

Supongamos que, debido a su complejidad, no podemos simular muestras de $p(x, y)$; pero que si es posible generar muestras de las condicionales $p(x|y)$ y $p(y|x)$. El muestreo de Gibbs consiste en, a partir de un x_0 , iterar entre:

$$y_k \sim p(y|x_k), \quad x_{k+1} \sim p(x|y_k)$$

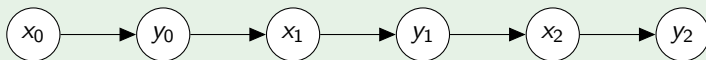
Luego de suficientes simulaciones (resultado asintótico), el último par (x, y) estará distribuido (aproximadamente) por $p(x, y)$.

Markov Chain Monte-Carlo (MCMC)



Técnicas de Muestreo

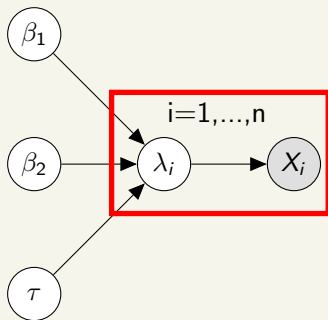
Markov Chain Monte-Carlo (MCMC)



Implementaciones más sofisticadas

- 1 Comenzar en la posición actual.
- 2 Proponer mudarse a una nueva posición cercana a la actual.
- 3 Aceptar/Rechazar la nueva posición basándose en el coherencia de la posición con los datos y distribuciones anteriores.
- 4
 - ▶ Si acepta: Pasar a la nueva posición. Regresar al Paso 1.
 - ▶ De lo contrario: no moverse de la posición actual. Regrese al Paso 1.
- 5 Después de una gran cantidad de iteraciones, se reportan todas las posiciones aceptadas.

PYMC: Programación Probabilística con Monte-Carlo

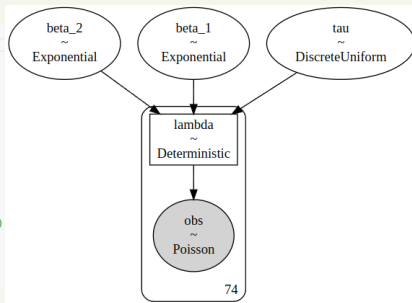



```
import pymc as pm
import numpy as np
import matplotlib.pyplot as plt

count_data = np.loadtxt("txtdata.csv")
n_count_data = len(count_data)

with pm.Model() as model:
    alpha = 1.0/count_data.mean()
    beta_1 = pm.Exponential("beta_1", alpha)
    beta_2 = pm.Exponential("beta_2", alpha)
    tau = pm.DiscreteUniform("tau", lower=0, upper=n_count_data - 1)
    idx = np.arange(n_count_data) # Index
    lambda_ = pm.Deterministic("lambda", pm.math.switch(tau > idx, beta_1, beta_2))
    observation = pm.Poisson("obs", lambda_, observed=count_data)
    trace = pm.sample(draws=1000, chains=2)

pm.model_to_graphviz(model)
```



PYMC

```
import pymc as pm
import numpy as np
import matplotlib.pyplot as plt

count_data = np.loadtxt("txtdata.csv")
n_count_data = len(count_data)

with pm.Model() as model:
    alpha = 1.0/count_data.mean()
    beta_1 = pm.Exponential("beta_1", alpha)
    beta_2 = pm.Exponential("beta_2", alpha)
    tau = pm.DiscreteUniform("tau", lower=0, upper=n_count_data - 1)
    idx = np.arange(n_count_data) # Index
    lambda_ = pm.Deterministic("lambda", pm.math.switch(tau > idx, beta_1, beta_2))
    observation = pm.Poisson("obs", lambda_, observed=count_data)
    trace = pm.sample(draws=1000, chains=2)

pm.model_to_graphviz(model)
```

```
beta_1_samples = trace.posterior['beta_1']
beta_2_samples = trace.posterior['beta_2']
tau_samples = trace.posterior['tau']
lambda_samples = trace.posterior['lambda']

with model:
    posterior_pred = pm.sample_posterior_predictive(trace)

pred_samples = posterior_pred.posterior_predictive['obs']

pm.plot_posterior(trace.posterior[['beta_1', 'beta_2', 'tau']], figsize=(7,4))
```

